

## مقایسه روش‌های جنگل تصادفی و رگرسیون لجستیک در پیش‌بینی مرگ‌ومیر مبتلایان به سرطان کولورکتال و عوامل مرتبط با آن

علی محمد کشت ورز حسام‌آبادی<sup>۱</sup>، ابراهیم حاجی‌زاده<sup>۲</sup>، محمدامین پورحسینقلی<sup>۳</sup>، احسان ناظم‌الحسینی مجرد<sup>۴</sup>

<sup>۱</sup>دانشجوی کارشناسی ارشد آمار زیستی، دانشکده علوم پزشکی، دانشگاه تربیت مدرس، تهران، ایران

<sup>۲</sup>استاد، گروه آمار زیستی، دانشکده علوم پزشکی، دانشگاه تربیت مدرس، تهران، ایران

<sup>۳</sup>دانشیار، مرکز تحقیقات گوارش و کبد، پژوهشکده بیماری‌های گوارش و کبد، دانشگاه علوم پزشکی شهید بهشتی، تهران، ایران

<sup>۴</sup>دکترای پزشکی مولکولی، مرکز تحقیقات گوارش و کبد، پژوهشکده بیماری‌های گوارش و کبد، دانشگاه علوم پزشکی شهید بهشتی، تهران، ایران

نویسنده رابط: ابراهیم حاجی‌زاده، نشانی: تهران، جلال آل احمد، پل نصر، دانشگاه تربیت مدرس، دانشکده علوم پزشکی، گروه آمار زیستی، تلفن: ۸۲۸۳۸۱۰.

پست الکترونیک: hajizadeh@modares.ac.ir

تاریخ دریافت: ۹۷/۰۳/۱۱؛ پذیرش: ۹۷/۰۷/۰۷

**مقدمه و اهداف:** هدف این مطالعه پیش‌بینی مرگ‌ومیر حاصل از سرطان کولورکتال در بیماران ایرانی و تعیین عوامل موثر بر آن با استفاده از روش جنگل تصادفی و رگرسیون لجستیک است.

**روش کار:** از اطلاعات ۳۰۴ بیمار مبتلا به سرطان کولورکتال ثبت مرکز تحقیقات گوارش و کبد دانشگاه علوم پزشکی شهید بهشتی طی سال‌های ۸۸ تا ۹۳ به صورت یک مطالعه گذشته‌نگر استفاده شد. تجزیه و تحلیل داده‌ها با استفاده از روش‌های جنگل تصادفی و رگرسیون لجستیک انجام شد. برای تحلیل داده‌ها از نرم افزار R نسخه ۳.۴.۳ استفاده شد.

**یافته‌ها:** ده متغیر مهمی که با مرگ‌ومیر سرطان کولورکتال ارتباط دارند، توسط روش جنگل تصادفی انتخاب شدند. چندین معیار مانند مساحت زیر منحنی مشخصه عملکرد (AUC) برای مقایسه روش جنگل تصادفی با رگرسیون لجستیک در نظر گرفته شد. با توجه به هر دو معیار، پنج متغیر اثرگذار رتبه‌بندی شده توسط جنگل تصادفی عبارتند از: مرحله سرطان، سن تشخیص بیماری، سن بیمار، مکانیسم فرار از سلول ایمنی و درجه تمایز یافتگی تومور. از نظر معیارهای مختلف، روش جنگل تصادفی عملکرد بهتری نسبت به رگرسیون لجستیک داشت (مساحت زیر منحنی ROC به ترتیب برای روش جنگل تصادفی و رگرسیون لجستیک برابر با ۰/۹۸، ۰/۸۰).

**نتیجه‌گیری:** متغیرهای سن تشخیص، مرحله سرطان، سن بیمار، مکانیسم فرار از سلول ایمنی و درجه تمایز یافتگی تومور به عنوان مهمترین عوامل موثر بر مرگ‌ومیر در سرطان کولورکتال به حساب می‌آیند که با تشخیص زودرس سرطان با برنامه‌های بیماریابی و غربالگری می‌توان بر طول عمر بیماران افزود.

**واژگان کلیدی:** سرطان کولورکتال، جنگل تصادفی، رگرسیون لجستیک، مرگ‌ومیر

### مقدمه

سرطان یکی از عوامل اصلی مرگ‌ومیر در جهان است که به‌عنوان یک مسئله مهم و قابل‌توجه در حوزه بهداشت عمومی شناخته می‌شود (۱). این بیماری دومین علت مرگ‌ومیر در کشورهای توسعه‌یافته و سومین عامل مرگ‌ومیر در کشورهای کمتر توسعه‌یافته است (۲).

سرطان در ایران پس از بیماری‌های قلبی-عروقی و سوانح، سومین عامل مرگ‌ومیر به‌شمار می‌آید. بر طبق گزارش‌های سالانه مرکز ثبت سرطان در ایران، سرطان کولورکتال پس از سرطان‌های معده، مثانه و پروستات در مردان، چهارمین سرطان است و در زنان بعد از سرطان سینه، دومین سرطان شایع

محسوب می‌شود (۳، ۴). در ایران هر سال سرطان کولورکتال، بروزی حدود ۶ تا ۷/۹ در ۱۰۰ هزار نفر جمعیت دارد (۵) و با میزان مرگی حدود ۱۱۹۸ در ۱۰۰ هزار نفر جمعیت تقریباً ۱۳٪ مرگ‌های ناشی از سرطان گوارش و ۵/۳٪ از مرگ‌ومیر ناشی از علل غیر حوادث در ایران را شامل می‌شود (۶).

سرطان کولورکتال یکی از انواع سرطان‌هاست که شامل سرطان روده بزرگ (سرطان کولون) و سرطان راست‌روده (سرطان رکتوم) و در اثر رشد کنترل نشده لایه داخلی اندام‌های کولون و رکتوم ایجاد می‌گردد که با افزایش سن افراد، افزایش می‌یابد و این

اطلاعات بیماران از سال ۱۳۸۸ تا ۱۳۹۳ ثبت شده بود. مرگ بیمار از طریق تماس با خانواده و اطرافیان بیمار مورد بررسی قرار گرفت. این طرح در کمیته اخلاق دانشگاه علوم پزشکی شهید بهشتی تصویب شده بود. برای همه بیماران ویژگی‌های بالینی و جمعیت شناختی شامل سن در زمان تشخیص بیماری، جنسیت، سابقه خانوادگی سرطان، سن بیمار، محل تومور، مرحله سرطان، درجه تمایزیافتگی تومور، مکانیسم فرار از سلول ایمنی، سابقه بیماری دیابت، سابقه بیماری آنمی، سابقه بیماری IBD، نتایج MSI، جهش BRAF، جهش KRAS، بتا کاتنین (B Catenin) و وضعیت حیاتی در تحلیل به کار رفت.

رگرسیون لجستیک<sup>۱</sup> یک ابزار تحلیلی بسیار عمومی است که اغلب در تحقیقات پزشکی برای پیش‌بینی متغیر پاسخ دوحالتی از طریق متغیرهای کمکی و عوامل، مورد استفاده قرار می‌گیرد (۱۴). برای یک متغیر پاسخ ۲، یک بردار از متغیرها  $X$  و  $\pi(X)$  که نشان‌دهنده شانس موفقیت با توجه به مقدار خاصی از  $X$  است، مدل رگرسیون لجستیک می‌تواند به صورت زیر نمایش داده شود.

$$\text{logit}[\pi(X)] = \log \left[ \frac{\pi(X)}{1 - \pi(X)} \right] = \alpha + \beta X$$

رگرسیون لجستیک از یک تبدیل لجیت برای محاسبه شانس یک متغیر پاسخ دوحالتی استفاده می‌کند (۱۵).

جنگل تصادفی<sup>۲</sup> یک روش یادگیری گروهی برای طبقه‌بندی است که توسط بریمن و همکاران در سال ۲۰۰۱ برای افزایش دقت طبقه‌بندی ارائه شد (۱۶). جنگل تصادفی یک طبقه‌بندی گروهی است که از تعدادی درخت تصمیم‌گیری تشکیل شده است که هر درخت با استفاده از یک نمونه بوت استرپ از داده‌های اصلی ساخته شده است. یک درخت به وسیله افزاز بازگشتی نمونه بوت استرپ بر اساس بهینه‌سازی یک قانون تقسیم‌بندی، رشد می‌کند (۱۷). تجزیه و تحلیل جنگل تصادفی شامل ۱۰۰۰ درخت است که هر کدام شامل ۴ متغیر ( $\sqrt{P}$ ) (که  $P$  تعداد متغیرهای مستقل است) می‌باشند (۱۸). از روش  $k$ -fold برای بررسی اعتبار نتایج و مقایسه آن با افراد در رگرسیون لجستیک، استفاده شد. مجموعه داده آزمون شامل یک سوم (۳۳ درصد) از کل داده‌ها و بقیه داده‌ها هم جزو مجموعه داده آموزشی می‌باشند. الگوریتم یادگیری را بر روی مجموعه داده آزمون اعمال می‌کنیم (به کار می‌بریم) و از مجموعه داده آموزشی برای الگوریتم تحت نظارت

سرطان با نام سرطان روده بزرگ شناخته می‌شود (۸، ۷). بر اساس گزارش‌های، میزان بروز سرطان کولورکتال در ایران طی ۲۵ سال اخیر روند رو به رشدی داشته و در مقایسه با کشورهای غربی جمعیت جوان‌تری از کشور را تحت تأثیر خود قرار داده است (۹)، در نتیجه اهمیت بررسی آن به عنوان یک مسئله بهداشت عمومی حائز اهمیت است. بررسی مرگ‌ومیر بیماران مبتلا به سرطان، یکی از انواع مطالعاتی است که وضعیت بیماری و عوامل مرتبط با آن را مشخص می‌کند.

در سال‌های اخیر توجه فراوانی به مدل‌های آماری برای طبقه‌بندی داده‌های پزشکی با توجه به بیماری‌های مختلف و پیامدهای آنها شده است. طبقه‌بندی یکی از مهم‌ترین کاربرد روش‌های آماری در علوم مختلف است که هدف پیش‌بینی یک پاسخ دوحالتی یا چندحالتی بر اساس برخی متغیرهای مستقل در مورد موضوعات مختلف است (۱۰).

در آمار عمده روش‌های مدل‌سازی تعیین روابط بین متغیرها، تعیین متغیرهای اثرگذار و پیش‌بینی است. به دلیل پیش‌فرض‌های موجود برای مدل‌های آماری کلاسیک و پیچیدگی تفسیر نتایج آنها برای محققان پزشکی در دهه‌های اخیر تکنیک‌های دیگری رواج یافته است. جنگل تصادفی یکی از روش‌های آماری نوین محسوب می‌شود که محدودیت‌های مدل‌های آماری کلاسیک را ندارد.

جنگل تصادفی یک روش ناپارامتری و متعلق به خانواده‌ی روش‌های گروهی است که عملکرد امیدوارکننده و نویدبخشی را در مطالعات مختلف نشان می‌دهد (۱۱). از ویژگی‌های مهم جنگل تصادفی عملکرد بالای آن در اندازه‌گیری اهمیت متغیرها برای مشخص کردن اینکه هر متغیر چه نقشی در پیش‌بینی پاسخ دارد، است (۱۲).

بنابراین با توجه به مطالب مطرح شده و افزایش بروز سرطان کولورکتال در چند دهه گذشته (۱۳)، مطالعه حاضر باهدف پیش‌بینی مرگ‌ومیر حاصل از سرطان کولورکتال با استفاده از روش جنگل تصادفی انجام شد. همچنین ما عملکرد جنگل تصادفی را با رگرسیون لجستیک مقایسه کردیم.

## روش کار

در این مطالعه گذشته‌نگر، پرونده ۳۰۴ بیمار مبتلا به سرطان کولورکتال که به بیمارستان طالقانی تهران مراجعه کرده بودند و اطلاعاتشان در پژوهشکده بیماری‌های گوارش و کبد دانشگاه علوم پزشکی شهید بهشتی ثبت شده بود، مورد بررسی قرار گرفت.

<sup>۱</sup> Logistic Regression

<sup>۲</sup> Random Forest

(۲۰).

### روش ارزیابی

برای مقایسه قدرت تمایز این دو مدل از مساحت زیر منحنی مشخصه عملکرد (AUC)، دقت (Accuracy)، ویژگی (Specificity) و حساسیت (Sensitivity) استفاده می‌کنیم. برای تجزیه و تحلیل داده‌ها از نرم‌افزار R استفاده شد. قبل از اینکه به این معیارها پرداخته شود لازم است به مفهوم ماتریس ابهام<sup>۶</sup> پرداخته شود. این ماتریس اجازه می‌دهد تجسمی از کارایی یک مدل ایجاد شود و حاوی اطلاعاتی در مورد دسته‌بندی واقعی و پیش‌بینی‌های انجام‌شده توسط یک سیستم دسته‌بندی است (۲۱). جدول شماره ۱ نشان‌دهنده ماتریس ابهام در یک مسئله دسته‌بندی دو کلاسه است.

جدول شماره ۱ - ماتریس ابهام

مجموع کل	داده‌های واقعی		پیش‌بینی
	منفی	مثبت	
TP+FP	FP	TP	مثبت
TN+FN	TN	FN	منفی
n	FP+TN	TP+FN	مجموع کل

در حوزه پزشکی منظور از مثبت، ابتلا به بیماری (در اینجا مرگ) و منظور از منفی، سالم (در اینجا زنده‌بودن) است.

▪ TP به معنی تعداد افرادی که مرده هستند و مدل پیش‌بینی کننده مرده پیش‌بینی کرده است.

▪ FP به معنی تعداد افرادی که زنده هستند و مدل پیش‌بینی کننده مرده پیش‌بینی کرده است.

▪ FN به معنی تعداد افرادی که مرده هستند و مدل پیش‌بینی کننده زنده پیش‌بینی کرده است.

▪ TN به معنی تعداد افرادی که زنده هستند و مدل پیش‌بینی کننده زنده پیش‌بینی کرده است.

با استفاده از مقادیر این ماتریس می‌توان مقدار معیارهای حساسیت<sup>۷</sup>، ویژگی<sup>۸</sup>، دقت<sup>۹</sup>، منفی کاذب<sup>۱۰</sup> و مثبت کاذب<sup>۱۱</sup> را

استفاده می‌کنیم. تجزیه و تحلیل جنگل تصادفی با استفاده از بسته "randomForest" در نرم‌افزار R انجام شد. یک RF آن‌قدر بزرگ است که تفسیر آن کار بسیار دشواری است، لذا نیازمند خلاصه کردن اطلاعات آن با استفاده از شاخص‌های کمی هستیم. یکی از این شاخص‌ها اهمیت متغیر<sup>۱</sup> (VI) است. VI شاخصی برای رتبه‌بندی متغیرها برحسب اهمیت آنها در اثرگذاری روی پاسخ است. معروف‌ترین شاخص‌های VI، شاخص اهمیت جینی<sup>۲</sup> و شاخص اهمیت جایگشتی<sup>۳</sup> است (۱۹).

شاخص اهمیت جایگشتی: برای محاسبه این شاخص، الگوریتم RF از تمام مشاهدات نمونه برای ساخت درخت استفاده نمی‌کند بلکه یک نمونه تصادفی با جایگذاری به حجم  $n_1$  (معمولاً برابر  $\frac{2}{3}n$ ) از مشاهدات انتخاب می‌شود. به مشاهدات انتخاب‌شده نمونه آموزشی<sup>۴</sup> (LS) و به بقیه آنها نمونه خارج کیسه<sup>۵</sup> (OOB) گفته می‌شود. درخت‌ها با مشاهدات LS ساخته می‌شوند و از OOB برای اندازه‌گیری ناخالصی درخت استفاده می‌شود. در هر درخت ابتدا اندازه ناخالصی روی مشاهدات OOB محاسبه می‌شود. سپس مقادیر متغیر  $X_i$  مشاهدات OOB به‌طور تصادفی جابجا می‌شوند و اندازه ناخالصی درخت روی مقادیر جابجا شده محاسبه می‌شود. اندازه اهمیت متغیر  $X_i$  در هر درخت، اختلاف بین این دو اندازه ناخالصی است و میانگین این مقادیر شاخص اهمیت جایگشتی است. انگیزه این روش این است که اگر  $X_i$  متغیر مهمی باشد جابجا شدن مقادیر آن به‌طور تصادفی منجر به افزایش ناخالصی درخت می‌شود در حالی که اگر متغیر تأثیرگذاری نباشد، تغییری در ناخالصی ایجاد نمی‌شود (۱۹).

از لحاظ تئوری (نظری)، جنگل تصادفی دارای چندین مزیت نسبت به رگرسیون است. مزیت اول، الگوریتم جنگل تصادفی می‌تواند متغیرهای مهم را به‌طور خودکار بدون توجه به اینکه چند متغیر در ابتدا استفاده شده، انتخاب کند که متفاوت از انتخاب مرحله‌ای (گام‌به‌گام) در رگرسیون لجستیک است. مزیت دوم، مقادیر گم‌شده و همچنین داده‌های نامتعادل می‌تواند به‌صورت خودکار توسط جنگل تصادفی به کار گرفته شود. مزیت سوم، جنگل تصادفی در مقایسه با رگرسیون لجستیک در مجموعه داده‌های بزرگ و با تعداد متغیر زیاد، عملکرد بهتری دارد

<sup>۶</sup> Confusion matrix

<sup>۷</sup> Sensitivity

<sup>۸</sup> Specificity

<sup>۹</sup> Accuracy

<sup>۱۰</sup> False Negative Rate

<sup>۱۱</sup> False Positive Rate

<sup>۱</sup> Variable Importance

<sup>۲</sup> Gini Importance Index

<sup>۳</sup> Permutation Importance Index

<sup>۴</sup> Learning Sample

<sup>۵</sup> Out-Of-Bag

محاسبه نمود. در ادامه به تعریف این معیارها می‌پردازیم.

### حساسیت

منظور از معیار حساسیت نسبت تعداد افراد مرده‌ای است که به درستی مرده دسته‌بندی شده‌اند به کل افرادی که مرده‌اند. با توجه به اینکه هرچقدر مقدار این معیار بزرگ‌تر باشد نشان‌دهنده دقت تشخیص افراد مرده است، مقدار این معیار در حوزه پزشکی بسیار بااهمیت است. مقدار این معیار با استفاده از رابطه زیر قابل محاسبه است.

$$\text{Sensitivity} = \frac{TP}{FN + TP}$$

### ویژگی

منظور از معیار ویژگی نسبت تعداد افراد زنده‌ای است که به درستی زنده دسته‌بندی شده‌اند به کل افراد زنده. این معیار با معیار حساسیت توازن<sup>۱</sup> دارند. معیار ویژگی از جمله معیارهایی است که در حوزه پزشکی دقت و کارایی مدل یادگیری را با استفاده از آن می‌سنجند و مقدار آن معیار با استفاده از رابطه زیر قابل محاسبه است.

$$\text{Specificity} = \frac{TN}{FP + TN}$$

### دقت

منظور از معیار دقت نسبت تعداد افراد مرده و زنده‌ای است که به درستی دسته‌بندی شده‌اند به کل جمعیت. با توجه به اینکه دو معیار حساسیت و ویژگی با یکدیگر توازن دارند از این معیار به‌عنوان معیار تخمین استفاده می‌شود. مقدار این معیار با استفاده از رابطه زیر قابل محاسبه است.

$$\text{Accuracy} = \frac{TN + TP}{FP + FN + TN + TP}$$

### منفی کاذب

منظور از معیار منفی کاذب نسبت تعداد افراد مرده‌ای است که زنده دسته‌بندی شده‌اند به کل افراد مرده. مقدار این معیار با استفاده از رابطه زیر قابل محاسبه است.

$$\text{FNR} = \frac{FN}{FN + TP}$$

### مثبت کاذب

منظور از معیار مثبت کاذب نسبت تعداد افراد زنده‌ای است که مرده دسته‌بندی شده‌اند به کل افراد زنده. مقدار این معیار با استفاده از رابطه زیر قابل محاسبه است.

$$\text{FPR} = \frac{FP}{FP + TN}$$

برای تجزیه و تحلیل داده‌ها از نرم‌افزار R نسخه ۳.۴.۳ استفاده شد. تجزیه و تحلیل جنگل تصادفی با استفاده از بسته "randomForest" انجام شد. همچنین از بسته "ROCR" برای رسم منحنی مشخصه عملکرد و محاسبه شاخص‌های نیکویی برازش استفاده کردیم.

### یافته‌ها

از تعداد ۳۰۴ بیمار مبتلا به سرطان کولورکتال ۱۵۶ نفر (۵۱/۳ درصد) مرد و ۱۴۸ نفر (۴۸/۷) زن بودند. میانگین سنی بیماران ۱۱/۱ ± ۵۶/۰۶ و حداقل و حداکثر سن به ترتیب ۲۸ و ۸۸ سال بود. در روش جنگل تصادفی بر اساس دو شاخص ضریب جینی و اهمیت جایگشتی، متغیرها بر پایه تأثیر آنها بر طبقه‌بندی مرتب شدند (۲۲). (شکل شماره ۱)

بر اساس شاخص ضریب جینی (سمت راست شکل ۱) سن در زمان تشخیص بیماری، سن بیمار، مرحله سرطان، درجه تمایز یافتگی تومور (Differentiation) و مکانیسم فرار از سلول ایمنی به ترتیب مهم‌ترین متغیرهای تأثیرگذار بر مرگ‌ومیر سرطان کولورکتال هستند. ترتیب مهم‌ترین متغیرهای تأثیرگذار بر اساس شاخص اهمیت جایگشتی (سمت چپ شکل ۱) هم به این صورت است: مرحله سرطان، سن بیمار، سن در زمان تشخیص بیماری، مکانیسم فرار از سلول ایمنی و درجه تمایز یافتگی تومور (Differentiation).

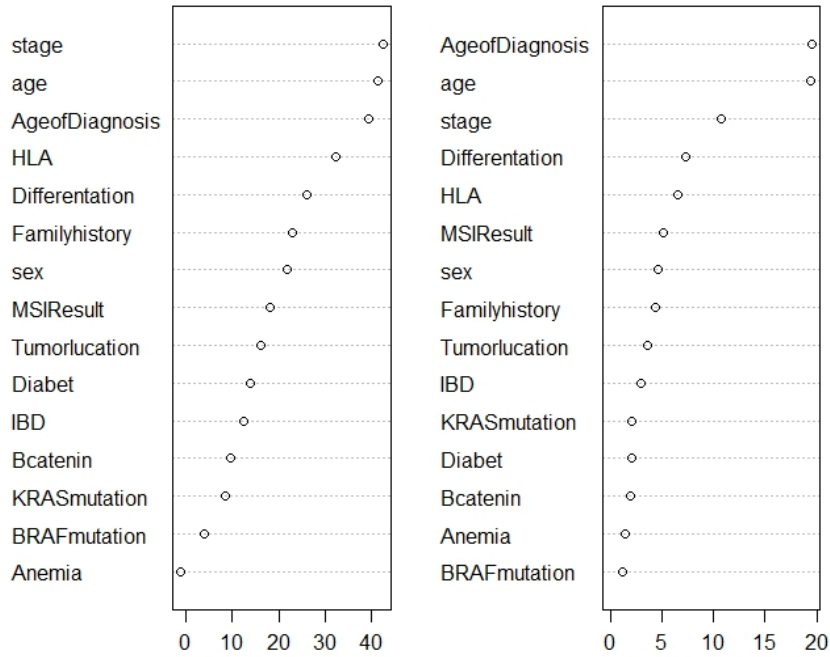
مدل رگرسیون لجستیک با استفاده از نرم‌افزار R شامل متغیرهای معنی‌دار حاصل از تحلیل تک متغیره، انجام شد. ورود متغیرها به مدل رگرسیون خطی تعمیم‌یافته بر اساس روش گام‌به‌گام (مرحله‌ای) و بهترین مدل با استفاده از معیار اطلاعات آکائیکی (AIC) انجام شد.

شکل شماره ۲ منحنی مشخصه عملکرد (ROC) برای روش‌های رگرسیون لجستیک و جنگل تصادفی را نشان می‌دهد. مساحت زیر منحنی مشخصه عملکرد برای روش‌های رگرسیون لجستیک و جنگل تصادفی به ترتیب برابر است با ۰.۸۰ و ۰.۹۸ که نشان‌دهنده

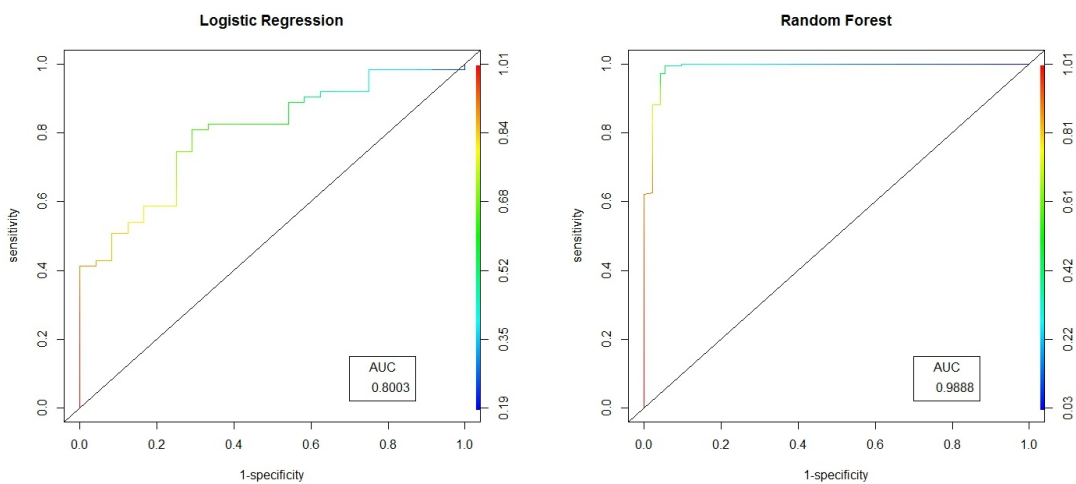
<sup>۱</sup> Trade off

نشان داده‌شده است که عملکرد بهتری برای روش جنگل تصادفی در مقایسه با روش رگرسیون لجستیک به دست آمد.

عملکرد بهتر روش جنگل تصادفی است اطلاعات در مورد ویژگی، حساسیت، دقت، درصد منفی کاذب و درصد مثبت کاذب برای روش‌های انجام‌شده در جدول شماره ۲



شکل شماره ۱ - شاخص ضریب جینی و شاخص اهمیت جایگشتی



شکل شماره ۲ - منحنی‌های مشخصه عملکرد برای روش‌های جنگل تصادفی و رگرسیون لجستیک

جدول شماره ۲ - دقت، ویژگی و حساسیت برای روش‌های رگرسیون لجستیک و جنگل تصادفی

شاخص‌های نیکویی برازش	جنگل تصادفی	رگرسیون لجستیک
حساسیت (Sensitivity)	۸۳٪	۶۰٪
ویژگی (Specificity)	۸۶٪	۷۱٪
دقت (Accuracy)	۸۵٪	۶۲٪
منفی کاذب (False Negative Rate)	۱۷٪	۴۰٪
مثبت کاذب (False Positive Rate)	۱۴٪	۲۹٪

## بحث

سنی زیر ۶۵ سال به‌طور معنی‌داری بیشتر از میانگین بقای افراد در گروه سنی بالای ۶۵ سال را گزارش داده است. همچنین در مطالعه ساکی مالچی و همکاران (۲۷) تأثیر سن در زمان تشخیص بر بقای بیماران تأیید شد. در مطالعاتی که توسط لیانگ و همکاران (۲۸) و چین و همکاران (۲۹) انجام شده است، سن در زمان تشخیص به‌عنوان عاملی مهم در مرگ بر اثر سرطان کولون گزارش شده است.

همچنین مطالعه‌ی ما نشان داد که مرحله بیماری بر خطر مرگ افراد در اثر سرطان کولورکتال، تأثیر معنی‌داری داشت. در مطالعه‌ی نصیری و همکاران (۲۶) مرحله بیماری عامل مهمی در پیش‌آگهی بیماران بود به‌طوری‌که بقا بیماران با بالا رفتن مرحله بیماری کاهش می‌یابد و هر چه مرحله‌ی بیماری پایین‌تر باشد، بقا بیماران بهتر خواهد بود. همچنین در مطالعه دهکردی و همکاران (۳۰) تأثیر مرحله‌ی بیماری بر بقا بیماران تأیید شد.

در مطالعه ما درجه تمایزیافتگی تومور اثر معنی‌داری بر مرگ افراد در اثر سرطان کولورکتال داشت. در مطالعه‌ای که توسط لیانگ و همکاران (۲۸) انجام شده است ارتباط معنی‌داری میان بقا بیماران مبتلا به سرطان کولون و درجه تمایزیافتگی مشاهده گردید. همچنین در مطالعات روشنایی و همکاران (۳۱)، آخوند و همکاران (۳۲) و دهکردی و همکاران (۳۳) تأثیر درجه تمایزیافتگی تومور بر بقا معنی‌دار بود.

از نقاط قابل توجه این مطالعه مقایسه دو روش آماری جنگل تصادفی و رگرسیون لجستیک است که می‌توان به توانایی بالاتر روش جنگل تصادفی نسبت به روش پرکاربرد رگرسیون لجستیک اشاره کرد و از محدودیت‌های این مطالعه می‌توان به عدم دسترسی به برخی اطلاعات بیماران اشاره کرد.

افزایش سرعت پیر شدن جوامع در کشورهای صنعتی، تعداد موارد سرطان از جمله سرطان کولورکتال را به‌سرعت افزایش می‌دهد. بروز سرطان کولورکتال در سه دهه اخیر روند صعودی را نیز در ایران نشان می‌دهد (۲۳). سرطان کولورکتال سالانه با تعداد نزدیک به یک‌میلیون مورد جدید و مرگ حدود ۵۰ درصد مبتلایان در پنج سال اول شروع بیماری به‌عنوان مشکل مهمی در سلامت عمومی شناخته می‌شود (۲۴). همین امر ضرورت مطالعه در مورد این سرطان را ایجاد می‌کند. در این مطالعه، بعضی از عوامل که ممکن است در مرگ‌ومیر بیماران سرطان کولورکتال مؤثر باشند از جمله سن در زمان تشخیص بیماری، جنسیت، سابقه خانوادگی سرطان، سن بیمار، محل تومور، مرحله سرطان، درجه تمایزیافتگی تومور، مکانیسم فرار از سلول ایمنی، سابقه بیماری آنمی، سابقه بیماری IBD، نتایج MSI، جهش BRAF، جهش KRAS، بتا کاتنین (B Catenin)، سابقه بیماری دیابت و وضعیت حیاتی مورد ارزیابی قرار گرفتند. از دو روش جنگل تصادفی و رگرسیون لجستیک استفاده کردیم. نتایج مطالعه ما نشان داد که از نظر معیارهای مختلف روش جنگل تصادفی دارای عملکرد بهتری نسبت به روش رگرسیون لجستیک بود. در مطالعه نوری و همکاران (۲۵) که از روش آماری ناپارامتری و نوین جنگل تصادفی استفاده کردند، بیان کردند که در تعیین متغیرهای مهمی که روی پاسخ تأثیرگذارند، روش جنگل تصادفی استنباط بهتری از متغیرها را می‌دهد.

با توجه به شکل شماره ۱، تجزیه و تحلیل اطلاعات نشان داد که متغیر سن در زمان تشخیص بیماری تأثیر بسزایی بر مرگ‌ومیر افراد مبتلا به سرطان کولورکتال داشت. مطالعه نصیری و همکاران (۲۶) رابطه معنی‌داری را بین سن در زمان تشخیص و بقای بیماران تأیید کرده‌اند به‌طوری‌که میانگین بقای افراد در گروه

## نتیجه‌گیری

دیگر، برخی از عوامل مؤثر بر مرگ‌ومیر بیماران مبتلا به سرطان کولورکتال به‌صورت مختلف و گوناگون گزارش شده است. لذا انجام مطالعات بیشتر جهت تعیین عوامل مؤثر بالینی و جمعیت شناختی بر مرگ‌ومیر بیماران مبتلا به سرطان کولورکتال می‌تواند مفید باشد. چون سن در زمان تشخیص و مرحله بیماری به‌عنوان عامل مؤثر در مرگ‌ومیر معرفی شد و مراحل بالای این بیماری می‌تواند با افزایش خطر مرگ همراه باشد، بنابراین آگاهی دادن به جامعه در جهت مراجعه هر چه سریع‌تر به پزشک و انجام معاینات، لازم و ضروری به نظر می‌رسد. همچنین با تشخیص زودرس سرطان با برنامه‌های بیماریابی و غربالگری مفید می‌توان بر طول عمر بیماران افزود.

در این مطالعه از دو روش جنگل تصادفی و رگرسیون لجستیک برای پیش‌بینی مرگ‌ومیر بیماران مبتلا به سرطان کولورکتال و تعیین عوامل خطر مرتبط با آن استفاده شد. بر اساس معیارهای حساسیت، ویژگی، دقت و مساحت زیر منحنی ROC روش جنگل تصادفی از توانایی بالاتری نسبت به رگرسیون لجستیک برخوردار بود؛ بنابراین روش جنگل تصادفی ابزار مناسب‌تری برای پیش‌بینی مرگ‌ومیر بیماران مبتلا به سرطان کولورکتال در این گونه داده‌ها است. باوجوداینکه مطالعات متعددی در حوزه‌ی سرطان کولورکتال در سراسر جهان انجام شده، ولی صحت عوامل تعیین‌شده هنوز هم جای بحث دارد. چون با توجه به مطالعات

## منابع

1. Azeem S, Gillani SW, Siddiqui A, Jandrajupalli SB, Poh V, Syed Sulaiman S. Diet and colorectal cancer risk in Asia-A systematic review. *Asian Pac J Cancer Prev.* 2015; 16: 5389-96.
2. Jemal A, Siegel R, Ward E, Hao Y, Xu J, Murray T, et al. Cancer statistics, 2008. *CA: a cancer journal for clinicians.* 2008; 58: 71-96.
3. Ahmadi A, Mosavi-Jarrahi A, Pourhoseingholi MA. Mortality determinants in colorectal cancer patients at different grades: a prospective, cohort study in Iran. *Asian Pacific journal of cancer prevention: APJCP.* 2015; 16: 1069-72.
4. Rezaianzadeh A, Safarpour AR, Marzban M, Mohaghegh A. A systematic review over the incidence of colorectal cancer in Iran. *Annals of colorectal research.* 2015; 3.
5. Sadjadi A, Malekzadeh R, Derakhshan MH, Sepehr A, Nouraei M, Sotoudeh M, et al. Cancer occurrence in Ardabil: Results of a population-based Cancer Registry from Iran. *International journal of cancer.* 2003; 107: 113-8.
6. Ganji A, Safavi M, Nouraei S, Nasser-Moghadam S, Merat S, Vahedi H, et al. Digestive and liver diseases statistics in several referral centers in Tehran, 2000-2004. *Govaresh.* 2006; 11: 3-8.
7. Baghestani AR, Daneshvar T, Pourhoseingholi MA, Asadzade H. Survival of colorectal cancer patients in the presence of competing-risk. *Asian Pac J Cancer Prev.* 2014; 15: 6253-5.
8. Li C, Lu H-J, Na F-F, Deng L, Xue J-X, Wang J-W, et al. Prognostic role of hypoxic inducible factor expression in non-small cell lung cancer: a meta-analysis. *Asian Pacific Journal of Cancer Prevention.* 2013; 14: 3607-12.
9. Bishehsari F, Mahdavinia M, Vacca M, Malekzadeh R, Mariani-Costantini R. Epidemiological transition of colorectal cancer in developing countries: environmental factors, molecular pathways, and opportunities for prevention. *World journal of gastroenterology: WJG.* 2014; 20: 6055.
10. Duda RO, Hart PE, Stork DG. *Pattern classification: John Wiley & Sons;* 2012.
11. Breiman L. *Random forests. Machine learning.* 2001; 45: 5-32.
12. Ishwaran H, Kogalur UB, Blackstone EH, Lauer MS. *Random survival forests. The annals of applied statistics.* 2008: 841-60.
13. Ansari R, Amjadi H, Norozbeigi N, Zamani F, Mir-Nasser S, Khaleghnejad A, et al. Survival analysis of colorectal cancer in patients underwent surgical operation in Shariati and Mehr Hospital-Tehran, in a retrospective study. *Govaresh.* 2007; 12: 7-15.
14. Harrell FE. *Ordinal logistic regression. Regression modeling strategies: Springer;* 2001, 331-43.
15. Hosmer Jr DW, Lemeshow S, Sturdivant RX. *Applied logistic regression: John Wiley & Sons;* 2013.
16. Svetnik V, Liaw A, Tong C, Culberson JC, Sheridan RP, Feuston BP. *Random forest: a classification and regression tool for compound classification and QSAR modeling. Journal of chemical information and computer sciences.* 2003; 43: 1947-58.
17. Tufféry S. *Data mining and statistics for decision making: Wiley Chichester;* 2011.
18. Strobl C, Boulesteix A-L, Kneib T, Augustin T, Zeileis A. *Conditional variable importance for random forests. BMC bioinformatics.* 2008; 9: 307.
19. Breiman L. *Classification and regression trees: Routledge;* 2017.
20. Geng M. *A comparison of logistic regression to random forests for exploring differences in risk factors associated with stage at diagnosis between black and white colon cancer patients: University of Pittsburgh;* 2006.
21. Kohavi R. *Glossary of terms. Special issue on applications of machine learning and the knowledge discovery process.* 1998; 30: 127-32.
22. Boulesteix A-L, Bender A, Lorenzo Bermejo J, Strobl C. *Random forest Gini importance favours SNPs with large minor allele frequency: impact, sources and recommendations. Briefings in Bioinformatics.* 2011; 13: 292-304.
23. Hosseini SV, Izadpanah A, Yarmohammadi H. *Epidemiological changes in colorectal cancer in Shiraz, Iran: 1980–2000. ANZ journal of surgery.* 2004; 74: 547-9.
24. Newton K, Newman W, Hill J. *Review of biomarkers in colorectal cancer. Colorectal disease.* 2012; 14: 3-17.
25. Noori s, Nourijelyani k, Mohammad k, Niknam M, Mahmoudi M, Andonian L, et al. *Random forests analysis: A modern statistical method for screening in high-dimensional studies and its application in a population-based genetic association study %J Journal of North Khorasan University of Medical Sciences.* 2012; 3: 93-101.
26. Nasiri S, Soroush A, Karamnezhad M, Mehrkhani F, Mosafa S, Hedayat A. *Prognostic Factors in the Survival Rate of Colorectal Cancer Patients after Surgery.* 2010.

27. Saki Malehi A, Hajizadeh E, Fatemi R. Evaluation of prognostic variables for classifying the survival in colorectal patients using the decision tree. *Iranian Journal of Epidemiology*. 2012; 8: 13-9.
28. Liang H, Wang XN, Wang BG, Pan Y, Liu N, Wang DC, et al. Prognostic factors of young patients with colon cancer after surgery. *World Journal of Gastroenterology: WJG*. 2006; 12: 1458.
29. Chin CC, Wang JY, Yeh CY, Kuo YH, Huang WS, Yeh CH. Metastatic lymph node ratio is a more precise predictor of prognosis than number of lymph node metastases in stage III colon cancer. *International journal of colorectal disease*. 2009; 24: 1297-302.
30. MOGHIMI DB, SAFAEI A, ZALI MR. Survival rates and prognostic factors in colorectal cancer patients. 2008.
31. Roshanaei G, Komijani A, Sadighi A, Faradmal J. Prediction of survival in patients with colorectal cancer referred to the Hamadan MRI center using of Weibull parameter model and determination of its risk factors during 2005-2013. 2014.
32. Akhoond M, Kazemnejad A, Hajizadeh E, Fatemi S, Motlagh A. Investigation of Influential Factors Affecting Survival Rate of Patients with Colorectal Cancer using Copula Function. *Iranian Journal of Epidemiology*. 2011; 6: 40-9.
33. Moghimi-Dehkordi B, Safaee A, Zali MR. Prognostic factors in 1,138 Iranian colorectal cancer patients. *International journal of colorectal disease*. 2008; 23: 683-8.



# Comparison of Random Forest and Logistic Regression Methods in Predicting Mortality in Colorectal Cancer Patients and Its Related Factors

Keshtvarz Hesam Abadi AM<sup>1</sup>, Hajizadeh E<sup>2</sup>, Pourhoseingholi MA<sup>3</sup>, Nazemalhossein Mojarad E<sup>4</sup>

1- Masters Student of Biostatistics, Faculty of Medical Sciences, Tarbiat Modares University, Tehran, Iran

2- Professor of of Biostatistics, Faculty of Medical Sciences, Tarbiat Modares University, Tehran, Iran

3- Gastroenterology and Liver Diseases Research Center, Research Institute for Gastroenterology and Liver Diseases, Shahid Beheshti University of Medical Sciences, Tehran, Iran

4- Gastroenterology and Liver Disease Research Center, Shahid Beheshti University of Medical Sciences, Tehran, Iran

**Corresponding author:** Hajizadeh E, hajizade@modares.ac.ir

(Received 1 June 2018; Accepted 29 September 2018)

**Background and Objectives:** The purpose of this study was to predict the mortality rate of colorectal cancer in Iranian patients and determine the effective factors on the mortality of patients with colorectal cancer using random forest and logistic regression methods.

**Methods:** Data from 304 patients with colorectal cancer registry from the Gastroenterology and Liver Research Center of Shahid Beheshti University of Medical Sciences during the years 2009 to 2014 were used as a retrospective study. Data analysis was performed using random forest and logistic regression methods. To analyze the data, R software version 3.4.3 was considered.

**Results:** Ten important variables related to colorectal cancer deaths were selected by random forest method. Several criteria such as the area under the characteristic curve (AUC) were used to compare the random forest method with logistic regression. According to both criteria, five important variables ranked by random forest were Cancer stage, age of diagnosis, patient's age, HLA, and degree of differentiation (tumor differentiation). In terms of different criteria, the random forest method had better performance than logistic regression (Area under the ROC curve for random forest and logistic regression methods was: 98%; 80% respectively).

**Conclusion:** Variables such as Cancer stage, age of diagnosis, patient's age, HLA, and degree of differentiation are considered as the most important factors affecting mortality in colorectal cancer, that the patients' longevity can be increased with the early diagnosis of cancer and screening programs.

**Keywords:** Colorectal cancer, Random forest, Logistic regression, Mortality